تعرین عملی شماره دو

مبانی بازیابی اطلاعات و جستجو وب

دانیال بیاتی ______دانیال بیاتی ____

محمد جواد کفایتی _____

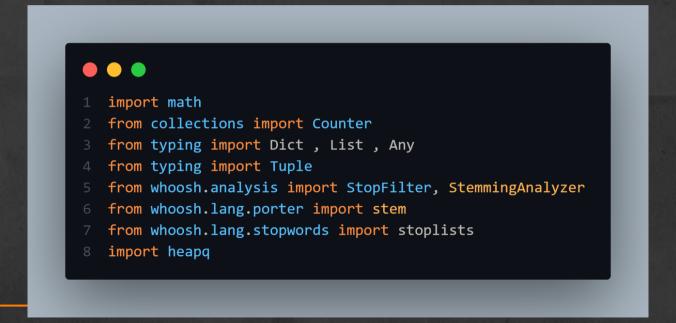
استاد: دکتر هدی مشایخی

تعدادی مجموعه داده مرجع برای بازیابی قرار دارد که توسط حل تمرین یکی از آنها به گروه شما تعلق میگیرد. در مجموعه داده، تعدادی سند، تعداد پرس و جو و همچنین ارتباط اسناد با پرس و جو ها مشخص شده است. برنامه ای بنویسید که بردار هر سند و هر پرس و جو را در فضای برداری با روش وزن دهی idf-Tf به دست آورد. سپس برای هر پرس و جو ۱۰ سند برتر مرتبط با پرس و جو را با روش شباهت کسینوسی محاسبه کنید و با مقایسه با مجموعه مرجع مقادیر صحت، یادآوری و معیار ۱ F ا را برای هر پرس و جو و میانگین برای کل مجموعه داده را محاسبه کنید



فایل <u>main.py</u> برنامه اصلی می باشد که از دو کلاس شباهت کسینوسی و جستوجو تشکیل شده است.

برای حل این تمرین از کتابخانه whoosh کمک گرفته شده است که توابع استفاده شده به صورت زیر می باشند:



```
1 class CosineSimilarity:
    def init (self, documents: List[Tuple[str, str]]) -> None:
        """ This function create list of vector for analyze content and scoring
        Args:
            documents (List[Tuple[str, str]]): List of document
        self.doc freqs : Dict[str, int] = Counter()
        self.create doc frequency(documents)
        self.idf: Dict[str, float] = {}
        self.create inverse term frequency(documents)
        self.doc vectors:Dict[int,Tuple[str,List[float]]] = {}
        self.create document vectors(documents)
    def create doc frequency(self, documents : List[Tuple[str, str]]) -> None:
        """ This function create list of document frequency that contain word and number of
            repeat in all documents
        Args:
            documents (List[Tuple[str, str]]): list of document
        for _, doc in documents:
            self.doc freqs.update(Counter(doc.split()))
    def create inverse term frequency(self, documents : List[Tuple[str, str]]) -> None:
        """ Create inverse term frequency with this formula
            idf for specified word = log(length of all documents / number of repeats in all docs)
        Args:
            documents (List[Tuple[str, str]]): List of document
        num docs = len(documents)
        for word, freq in self.doc freqs.items():
            self.idf[word] = math.log(num docs / freq)
```

: CosineSimilarity کلاس

تابع create_doc_frequency به عنوان ورودی لیست داکیومنت های مورد نظر را گرفته و ft هر کلمه را محاسبه کرده و آرایه doc_freq را بروزرسانی می کند.

ر create_inverse_term_frequency به عنوان ورودی لیستی از داکیومنت ها را دریافت کرده و با پیمایش درون آرایه docs_freqs فرکلمه را حساب کرده و ذخیره می کند.

```
def create document vectors(self, documents: List[Tuple[str, str]]) -> None:
        This functions create list of tuples that contain title and
        inverse term frequency each word * term frequency of each word
    Args:
        documents (List[Tuple[str, str]]): List of document
    for i, (title, doc) in enumerate(documents):
        tf = Counter(doc.split())
        vec = [tf[w] * self.idf[w] for w in self.doc freqs.kevs()]
        self.doc vectors[i] = (title, vec)
def cosine similarity(self, query :str ="query") -> List[Tuple[str, float]]:
     """ This func create list of the approximate amount query with each document
    Args:
        guery (str. optional): A string that the user enters to measure the similarity
    Returns:
        List[Tuple[str, float]]: List of the approximate amount query with each document
    # Removed stop words and finds the root of the rest
    tf: Dict[str, int] = Counter(query.split())
    q_vec:List[float] = [tf[w] * self.idf[w] if w in tf else 0 for w in self.doc_freqs.keys()]
    sims:List[Tuple[str, float]] = []
    for _, (title, d_vec) in self.doc_vectors.items():
        dot product:float = sum([a*b for a,b in zip(d vec, q vec)])
        norm d:float = math.sqrt(sum([a*a for a in d vec]))
        norm q:float = math.sqrt(sum([a*a for a in q vec]))
            sim:float = dot product / (norm d * norm q)
        except ZeroDivisionError:
            sim:float = dot_product / ((norm_d * norm_q)+0.0001)
        sims.append((title, sim))
    return sims
```

تابع tf-idf های حساب شده در دو تابع قبل توجه به tf-idf های حساب شده در دو تابع قبل بردار داکیومنت ها را می سازد تا در مرحله بعدی شباهت کسینوسی محاسبه شود.

تابع cosine_similarity پس از محاسبه بردار کوئری مورد نظر با محاسبه ضرب داخلی بردار ها و اندازه هر بردار شباهت داکیومنت ها به کوئری مورد نظر رامحاسبه کرده و درون آرایه sims ذخیره می کند تا شباهت همه داکیومنت ها به کوئری را داشته باشیم.

```
class Searcher:
   def init (self .address of documents:str = "./docs/MED.ALL"):
        self.address of documents = address of documents
       self.read content from docs()
       self.normalize contents()
        self.cosineSimilarity = CosineSimilarity(self.information)
        self.results: List[Tuple[str, int]] = []
    def preprocessed text(self, text:str )-> str:
        """This function takes a Text and removes the stop words and finds the root of the rest
       Args:
            text (str): A string of words in a row
       Returns:
            str: A string that removed stop words and finds the root of the rest
        stopwords: Any = frozenset(stoplists["en"])
        analyzer = StemmingAnalyzer(stemfn=stem, stoplist=stopwords) | StopFilter(stoplist=stopwords)
       # Tokenize and analyze the text using the defined analyzer
       tokens:List[str] = [token.text for token in analyzer(text)]
        # Join the tokens back together to form a preprocessed string
       preprocessed text = " ".join(tokens)
       return preprocessed text
    def read content from docs(self, address:str = "./docs/MED.ALL"):
        """ This func create list of information of each document that each info includes
           a document number and content
       Args:
            address (str. optional): " Address of file that content documents". Defaults to "./docs/MED.ALL".
       with open(address, 'r') as file:
           documents = file.read()
           list of documents = documents.split(".I")[1:]
           # Create list of document number and content
           self.information:Any = [[
                           doc.split(".W")[0].replace("\n","").strip(),
                           doc.split(".W")[1].replace("\n"," ")]
                           for doc in list_of_documents
```

کلاس Searcher:

تابع اول همانطور که در کد مشاهده می کنید وظیفه نرمال سازی داکیومنت ها را دارد لذا با استفاده از یه سری عملیات همچون stemming و ... داکیومنت های مورد نظر را نرمال می کنیم.

تابع دوم وظیفه خواندن از فایل را برعهده دارد که تمام داکیومنت های موجود در پوشه docs و فایل MED.ALL فایل و با فرمت قابل پیمایش ارائه میدهد

```
def normalize contents(self) :
        """ Normalize content of each document with call preprocessed text func
        Args:
           information (List[List[str]]): List of info that each element is a list that contain a document number and content
        Returns:
           List[List[str]]: A list that removed stop words and finds the root of the rest
       for index,info in enumerate(self.information):
           document number , content = info[0],info[1]
           # Update content each info with remove stop words and ...
           self.information[index] = (document number, self.preprocessed text(content))
    def search query(self, query:str)-> List[Any]:
        """ This function is to get similar document to a document
        Args:
           query (str): A text to get similar document
        Returns:
           List[tuple[str,int]]: results of similar document
       query = self.preprocessed text(query)
        self.results = self.cosineSimilarity.cosine similarity(query)
    def get nlargest similarity doc(self, n first:int) ->List[Tuple[str,int]]:
        return heapq.nlargest(n first, self.results, key=lambda x: x[1])
    def print results(self, n first = 10)-> None:
       if self.results != None :
           n first similarty doc = self.get nlargest similarity doc(n first)
           print('-----')
           print("
                       SCORE
                                            DOCUMENT NUMBER")
           print('-----')
           for index,result in enumerate(n first similarty doc):
               document number = result[0]
               score = index + 1
               print(f"
                            {score}
                                                      {document number}")
       else :
           print("Please First call search query function")
```

پس از نرمال سازی داکیومنت های موجود شباهت کسینوسی هر داکیومنت را با کوئری مورد نظر حساب کرده و درون آرایه result ذخیره می کنیم سپس با توجه به شباهت بیشتر ۱۰ تای بیشتر از نظر شباهت را جدا کرده و در خروجی نمایش می دهیم

```
Please enter your query: the crystalline lens in vertebrates, including humans.
              RESULTS --
  SCORE
                ***
                        DOCUMENT NUMBER
                        72
                        500
                                                             965
                        360
                                                               s = Searcher()
                        181
                                                                query = input('Please enter your query: ')
                        171
                                                                s.search_query(query)
                        166
                                                               s.print results(12)
  8
                        15
                        513
                         838
```

تست و بررسی داکیومنت ها :

برای مثال اگر کوئری زیر را داشته باشیم:

query: the crystalline lens in vertebrates, including humans.

خروجی داکیومنت های مورد نظر نشان می دهد که با بررسی فایل MED.REL صحت شباهت های به دست آمده به درستی برقرار است.